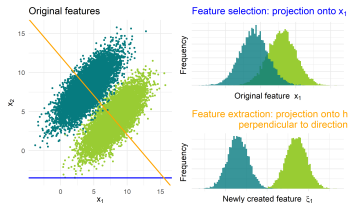


Feature Selection



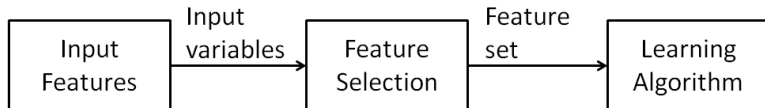
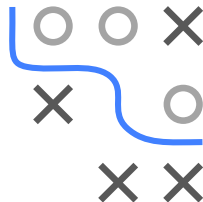
- Adding more features can be detrimental to predictive performance.
- Benefits of keeping only informative features for the model



INTRODUCTION

Feature selection deals with

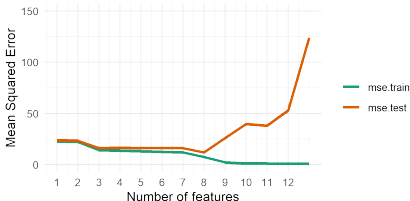
- techniques for choosing a suitable subset of features
- evaluating the influence of features on the model



Feature selection can be performed relying on domain knowledge and expert input, or using a data-driven algorithmic approach.

MOTIVATION

-
- A 3x3 grid of symbols. The top row contains 'o', 'o', 'x'. The middle row contains 'x', an empty space, 'o'. The bottom row contains an empty space, 'x', 'x'. A blue line starts at the top-left corner, goes right, then down, then right, separating the 'o's from the 'x's.



MOTIVATION

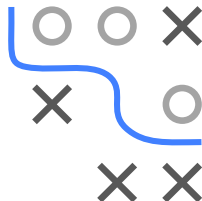
- In high-dimensional data sets, we often have prior information that many features are either irrelevant or redundant.
- Feature selection is critical for
 - reducing noise and overfitting,
 - improving performance/generalization,
 - interpretability by identifying most informative features.
- Feature selection can also remedy problems arising in small n regimes or under limited computational resources.
- Many models require $n > p$ data. Thus, we either need to
 - adapt models to high-dimensional data (e.g. regularization),
 - design entirely new procedures for $p > n$ data, or
 - use the preprocessing methods addressed in this lecture.



SIZE OF DATASETS

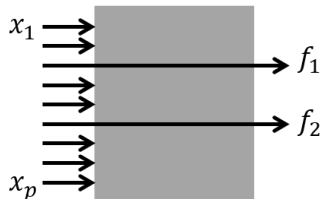
The increasingly automatized collection of information makes data sets with extremely high dimensionality available, while classical models were developed for small p data.

- **Classical setting:** Up to around 10^2 features, feature selection might be relevant, but benefits often negligible.
- **Datasets of medium to high dimensionality:** At around 10^2 to 10^3 features, classical approaches can still work well, while principled feature selection helps in many cases.
- **High-dimensional data:** 10^3 to 10^9 or more features. Examples are e.g. micro-array / gene expression data and text categorization (bag-of-words features). If, in addition, observations are few, the scenario is called $p \gg n$.



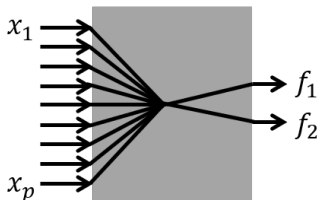
FEATURE SELECTION VS. EXTRACTION

Feature selection



- Creates a subset of original features \mathbf{x} by selecting $\tilde{p} < p$ features \mathbf{f} .
- Retains information on selected individual features.

Feature extraction



- Maps p features in \mathbf{x} to \tilde{p} extracted features \mathbf{f} .
- Info on individual features can be lost through (non-)linear combination.



A 3x3 grid with a blue path starting at the top-left cell (0,0) and ending at the bottom-right cell (2,2). The path moves right to (0,1), down to (1,1), right to (1,2), and down to (2,2). The cells (0,2), (1,0), and (2,0) are marked with 'X', while the other cells are empty.

-